# Video Representation Via 3D Shaped Mosaics

Pedro M. Q. Aguiar *
ISR, Instituto Superior Técnico
Lisboa, Portugal
aguiar@isr.ist.utl.pt

José M. F. Moura
ECE Dep, Carnegie Mellon University
Pittsburgh, PA
moura@ece.cmu.edu

## Abstract

*We generalize to 3D shaped mosaics the **Generative Video** representation of video sequences introduced by Jasinschi and Moura. Using a parametric representation of the 3D shape, we recover the 3D shape and 3D motions from the 2D motions in the video sequence. In this paper we consider piecewise planar object shapes under orthography and demonstrate our approach with a real life video clip.*

## 1 Introduction

**Generative Video (GV)**, introduced by Jasinschi and Moura, e.g., [1], reduces video sequences to world images and ancillary data. The world images represent the background and any moving objects, while the ancillary data describes for example the motions (camera and objects). In the original formulation of GV, the world images are modeled as simple planar scenarios. This representation fails when the relative depth of the scene structure is not negligible. In this paper we recover 3D world image representations for the video sequence. Within our framework, the major task is to recover the 3D structure (3D shape and 3D motion) from a 2D video sequence.

**Previous related work** Consider the case of a single rigid body object, moving relative to the camera, at a large distance when compared to the object depth. In this scenario, to recover the 3D shape by estimating the absolute depth is inaccurate. Tomasi and Kanade, e.g., [2], introduced a method to recover structure from motion without computing the absolute depth as an intermediate step. In the approach of Tomasi and Kanade, the object shape is represented by the 3D position of a set of feature points. The 2D projection of each feature point is tracked along the image sequence. The 3D shape and motion are then estimated by factorizing a measurement matrix whose entries are the set of trajectories of the feature point projections.

**Proposed approach** We represent parametrically the 3D shape of the rigid body object and apply *Maximum Likelihood* (ML) estimation. The observations are the orthographic projection of the object texture plus noise. The problem is now to estimate from the given sequence of images all the unknowns involved (3D shape parameters, 3D motion parameters, and object texture). We obtain a feasible approximation

to the ML estimator by showing that the parametric representation of the 3D rigid body shape induces a parametric model for the optical flow in the 2D image plane. We estimate the optical flow parameters using known techniques and apply *Least Squares* (LS) to resolve the 3D shape and motion parameters from the optical flow parameter estimates. Our technique was introduced in reference [3] where we considered that the world is 2D and the images are 1D projections of the world.

Our method relates to that of Tomasi and Kanade [2] in two very fundamental ways: we estimate directly the 3D shape, instead of computing depth as an intermediate step; and our algorithm leads to the factorization of a measurement matrix. Aside from this, our approach is different. We do not rely on the tracking of feature points. Instead, we use a parametric description of the 3D shape and recover the 3D structure from the parameterization induced in the optical flow. The advantage of our approach is two-fold. First, the tracking of feature points may be unreliable when processing noisy video sequences. The work in reference [2] assumes a very short interval between frames for easy feature tracking. We make no such assumption because large displacements are taken care of by a multiresolution approach to the estimation of the optical flow parameters. Second, we estimate 3D shape and motion from a sequence of few flow parameters instead of needing to process a large set of feature trajectories.

**Paper overview** Section 2 formulates the problem, derives the ML-based cost function, and solves for the texture estimate. Section 3 describes our approach to the recovery of the 3D shape and 3D motion for piecewise planar shapes. We demonstrate our approach by analyzing a real life video clip in section 4. Section 5 concludes the paper.

## 2 Problem Formulation

**Observation model** The frame $I_f$ captured at time $f$, is modeled as a noisy observation of the orthogonal projection $\mathcal{P}$ of the rigid object $\mathcal{O}$ (assumed segmented)

$$I_f = \mathcal{P}(\mathcal{O}, m_f) + W_f \tag{1}$$

$m_f$ defines the position and orientation of the rigid object relative to the camera coordinate system. For simplicity, the noise $W_f$ is white, Gaussian.

---

*The first author was partially supported by INVOTAN.

823

The object $\mathcal{O}$ is described by its 3D shape $\mathcal{S}$ and texture $\mathcal{T}$. We model the shape $\mathcal{S}$ by a parametric description $\mathcal{S}(a)$ of the surface of the object, where $a$ is an unknown vector. The texture $\mathcal{T}$ represents the light received by the camera after reflecting on the object surface. Texture depends on the object surface photometric properties, as well as on the environment illumination conditions. We assume the texture at a given surface point does not change with time. The operator $\mathcal{P}$ returns a real valued function defined over the image plane. $\mathcal{P}$ is a nonlinear mapping of $\mathcal{T}$ that depends on the object shape $\mathcal{S}$ and the object position $m_f$. The intensity level of the projection of the object at pixel $u$ on the image plane in terms of $\mathcal{T}$ is

$$\mathcal{P}\left(\mathcal{O}, m_f\right)(u) = \mathcal{T}\left(s_f(\mathcal{S}, m_f; u)\right) \qquad (2)$$

where $s_f(\mathcal{S}, m_f; u)$ is the nonlinear mapping that lifts the point $u$ on image $I_f$ to the corresponding point on the 3D object surface. This mapping $s_f(\mathcal{S}, m_f; u)$ is determined by the object shape $\mathcal{S}(a)$, and the position $m_f$. To simplify the notation, we will write explicitly only the dependence on $f$, i.e., $s_f(u)$. Let $u_f(s)$ be the inverse map of $s_f(u)$. The point $s$ on the surface of the object projects onto $u_f(s)$ on the image. The mapping $u_f(s)$, seen as a function of the frame index $f$, for a particular surface point $s$, is the trajectory of the projection of that point in the image, i.e., the motion induced in the image, usually referred to as *optical flow*. In the sequel, we refer to the map $u_f(s)$ as the optical flow map. The observation model (1) is rewritten by using (2) as

$$I_f = \mathcal{T}\left(s_f(u)\right) + W_f \qquad (3)$$

**ML estimate** Given the observation model (3), the 3D shape and the 3D motion of the object $\mathcal{O}$ are recovered from the video sequence $\{I_f, 1 \le f \le F\}$ by estimating all the unknown parameters: the 3D shape parameter $a$; the texture $\mathcal{T}$; and the set of 3D positions of the object $\{m_f, 1 \le f \le F\}$ with respect to the camera.

With the noise $W_f$ zero mean, white Gaussian, the ML estimate minimizes the cost function $C_{\mathrm{ML}}$ defined as

$$C_{\mathrm{ML}}\left(a, \mathcal{T}, \{m_f\}\right) = \sum_{f=1}^{F} \int \left[I_f(u) - \mathcal{T}\left(s_f(u)\right)\right]^2 du \qquad (4)$$

**Texture estimate** We show in [4] that the ML estimate $\widehat{\mathcal{T}}(s)$ that minimizes $C_{\mathrm{ML}}$ is

$$\widehat{\mathcal{T}}(s) = \frac{\sum_{f=1}^{F} I_f(u_f(s)) J_f(s)}{\sum_{f=1}^{F} J_f(s)} \qquad (5)$$

where the function $J_f(s)$ is the Jacobian of the mapping $u_f(s)$, $J_f(s) = |\nabla u_f(s)|$. Expression (5) states that the estimate of the texture of the object at the surface point $s$ is a weighted average of the measures of the intensity level corresponding to that surface point.

By inserting the texture estimate $\widehat{\mathcal{T}}$ in (4), we express $C_{\mathrm{ML}}$ in terms of the optical flow mappings $\{u_f(s)\}$. After manipulations, see [4], $C_{\mathrm{ML}}$ is written as

$$C_{\mathrm{ML}} = \sum_{f=2}^{F} \sum_{g=1}^{f-1} \int \left[I_f(u_f(s)) - I_g(u_g(s))\right]^2 \frac{J_f(s) J_g(s)}{\sum_{h=1}^{F} J_h(s)} ds \qquad (6)$$

Eliminating the dependence on the texture, we are left with a cost function that depends on the *structure* (3D shape $\mathcal{S}(a)$ and 3D motion $\{m_f\}$) only through the *motion* induced in the image plane, i.e., through the optical flow mappings $\{u_f(s)\}$. Recall that $u_f$ depends on the shape $\mathcal{S}$ and the motion $m_f$.

**Summary of the approach** To recover the 3D shape and the 3D motion of the object $\mathcal{O}$ from the image sequence, we do not attempt the direct minimization of $C_{\mathrm{ML}}$ over the parameters $a$ and $\{m_f\}$. Rather, we exploit the constraints induced on the optical flow by the orthogonality of the projection $\mathcal{P}$, the rigidity of the motions (rigid object), and the parameterization of the surface shape of the object. The constraints induced on the optical flow enable us to parameterize the optical flow mapping $u_f(s)$ in terms of a parameter vector $\alpha_f$ as $u(\alpha_f; s), 1 \le f \le F$. The parameter vector $\alpha_f$ is directly related to the 3D shape parameter $a$ and the 3D position $m_f$, as will be shown in section 3, i.e., $\alpha_f = \alpha(a, m_f)$. The steps of our approach to recover the 3D structure, i.e., the 3D shape parameter $a$ and the set of 3D positions $\{m_f\}$, are summarized as: **i)** Given the image sequence $\{I_f, 1 \le f \le F\}$, estimate the set of time varying optical flow vectors $\{\alpha_f, 1 \le f \le F\}$ parameterizing the optical flow mappings $\{u(\alpha_f; s)\}$. **ii)** From the sequence of estimates $\{\widehat{\alpha}_f, 1 \le f \le F\}$, solve for the shape parameter vector $a$ and the object positions $\{m_f, 1 \le f \le F\}$.

## 3 Piecewise Planar Surface

Attach a coordinate system to the object given by the axes labeled by $x$, $y$, and $z$. We consider objects whose shape is given by a piecewise planar surface with $K$ patches. The shape parameter vector $a$ collects the coefficients $a = \{a_{00}^k, a_{10}^k, a_{01}^k, 1 \le k \le K\}$ where

$$1 \le k \le K; \qquad z = a_{00}^k + a_{10}^k x + a_{01}^k y \qquad (7)$$

describes the shape of the patch $k$.

To capture the 3D motion of the rigid object, we attach a coordinate system to the camera given by the axes $u$, $v$, and $w$, see figure 1. We express the object position at time instant $f$ in terms of $\left(t_{uf}, t_{vf}, t_{wf}, \theta_f, \phi_f, \psi_f\right)$ where $\left(t_{uf}, t_{vf}, t_{wf}\right)$ are the coordinates of the origin of the object coordinate system with respect to the camera coordinate system (3D translation), and $(\theta_f, \phi_f, \psi_f)$ determine the orientation of the object coordinate system relative to the camera coordinate system (3D rotation).

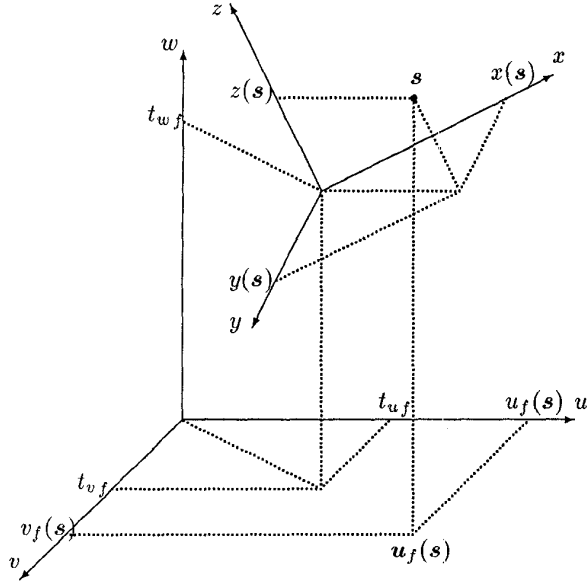The point $(x, y, z)$ projects at time instant $f$ on the

Figure 1: Object and camera coordinate systems.

image coordinates $u_f = [u_f, v_f]^T$ given by

$$u_f = \begin{bmatrix} i_{xf} & i_{yf} & i_{zf} \\ j_{xf} & j_{yf} & j_{zf} \end{bmatrix} \begin{bmatrix} x \\ y \\ z \end{bmatrix} + \begin{bmatrix} t_{uf} \\ t_{vf} \end{bmatrix} \qquad (8)$$

where the matrix that multiplies $[x, y, z]^T$ is a submatrix of the well known 3D rotation matrix, determined by the angles $(\theta_f, \phi_f, \psi_f)$:

$$i_{xf} = \cos\theta_f \cos\phi_f \qquad (9)$$

$$i_{yf} = \sin\phi_f \qquad (10)$$

$$i_{zf} = -\sin\theta_f \cos\phi_f \qquad (11)$$

$$j_{xf} = \sin\theta_f \sin\psi_f - \cos\theta_f \sin\phi_f \cos\psi_f \qquad (12)$$

$$j_{yf} = \cos\phi_f \cos\psi_f \qquad (13)$$

$$j_{zf} = \sin\theta_f \sin\phi_f \cos\psi_f + \cos\theta_f \sin\psi_f \qquad (14)$$

Expression (8) shows that the orthogonal projection is insensitive to the translation along the $w$ axis. This reflects the fact that under orthography the absolute depth can not be estimated. Only the set of positions $\{m_f = \{t_{uf}, t_{vf}, \theta_f, \phi_f, \psi_f\}, 1 \leq f \leq F\}$ can be estimated from the image sequence.

### 3.1 Optical flow

We show that the optical flow mappings $\{u_f(s)\}$ are described parametrically. Choose the coordinate $s = [s, r]^T$ of the texture function to coincide with the coordinates $[x, y]$ of the object coordinate system. Also, the object coordinate system and the

camera coordinate system coincide in the first frame. A point with coordinate $s = [s, r]^T$ in the object surface projects on $u = [x, y]^T = [s, r]^T = s$ in the first frame, so $u_1(s) = s$. At time $f$, that point projects according to (8). For a point $s$ that falls in patch $k$ of the object surface, we have

$$u_f(s) = \begin{bmatrix} i_{xf} & i_{yf} & i_{zf} \\ j_{xf} & j_{yf} & j_{zf} \end{bmatrix} \begin{bmatrix} s \\ r \\ a_{00}^k + a_{10}^k s + a_{01}^k r \end{bmatrix}$$

$$+ \begin{bmatrix} t_{uf} \\ t_{vf} \end{bmatrix} \qquad (15)$$

Define the set of parameters $\alpha_f^k = \{\alpha_{fmn}^{uk}, \alpha_{fmn}^{vk}\}$ as the coefficients of the powers of $s$ and $r$ above,

$$\alpha_{f00}^{uk} = t_{uf} + i_{zf} a_{00}^k, \qquad \alpha_{f00}^{vk} = t_{vf} + j_{zf} a_{00}^k \quad (16)$$

$$\alpha_{f10}^{uk} = i_{xf} + i_{zf} a_{10}^k, \qquad \alpha_{f10}^{vk} = j_{xf} + j_{zf} a_{10}^k \quad (17)$$

$$\alpha_{f01}^{uk} = i_{yf} + i_{zf} a_{01}^k, \qquad \alpha_{f01}^{vk} = j_{yf} + j_{zf} a_{01}^k \quad (18)$$

The optical flow between frames $I_1$ and $I_f$ in the image of the surface patch $k$ is written in terms of the optical flow parameter vector $\alpha_f^k$ as

$$u_f(s) = u(\alpha_f^k; s, r) = \begin{bmatrix} \alpha_{f00}^{uk} + \alpha_{f10}^{uk} s + \alpha_{f01}^{uk} r \\ \alpha_{f00}^{vk} + \alpha_{f10}^{vk} s + \alpha_{f01}^{vk} r \end{bmatrix}$$

The optical flow parametrization above is usually referred to as *affine motion model*. The ML estimation of $\{\alpha_f^k\}$ leads to the minimization of $C_{ML}$ given by (6) with respect to the set of vectors $\{\alpha_f^k, 1 \leq f \leq F, 1 \leq k \leq K\}$ parameterizing the mappings $\{u_f(s) = u(\alpha_f; s)\}$. In practice, this is a highly complex task. A more feasible and practical solution decouples the estimation of each vector $\alpha_f$ from the estimation of the remaining vectors $\alpha_g, g \neq f$, by simplifying the cost function (6). Instead of using all possible pairs of frames, compare all frames with respect to the first frame $I_1$. Also, neglect the weighting term $\frac{J_f(s)J_g(s)}{\sum_{h=1}^{F} J_h(s)}$, and obtain the usual expression for the optical flow estimation:

$$\{\widehat{\alpha_f^k}\} = \arg\min_{\{\alpha_f^k\}} \sum_{k=1}^{K} \int_{\mathcal{R}_k} \left[I_f\left(u(\alpha_f^k; s)\right) - I_1(s)\right]^2 ds$$

We compute the optical flow parameter vector estimates $\{\widehat{\alpha_f^k}\}$, by using known techniques, see references [5] and [6].

### 3.2 3D Structure from 2D Optical Flow

The set of equations (16-18) defines an overconstrained system with respect to the 3D shape parameters $\{a_{00}^k, a_{01}^k, a_{10}^k, 1 \leq k \leq K\}$ and to the 3D positions $\{t_{uf}, t_{vf}, \theta_f, \phi_f, \psi_f, 1 \leq f \leq F\}$. The estimate $\{\hat{a}_{mn}^k\}$ of the object shape and the estimate

$\left\{\hat{t}_{uf}, \hat{t}_{vf}, \hat{\theta}_f, \hat{\phi}_f, \hat{\psi}_f\right\}$ of the object positions are the *Least Squares* (LS) solution of the system. We first solve for the translation, leading to a closed-form solution. Then, replace the translation estimate and solve for the remaining motion parameters and shape parameters by using a two-step iterative method.

**Translation estimation** The translation components along the camera plane at instant $f$, $t_{uf}$ and $t_{vf}$ only affect, respectively, the set of parameters $\left\{\alpha_{f00}^{uk}\right\}$ and $\left\{\alpha_{f00}^{vk}\right\}$. If the parameters $\{a_{00}^k\}$ and $\{\theta_f, \phi_f, \psi_f\}$ are known, the LS estimate of $\{t_{uf}, t_{uf}\}$ is given by

$$\hat{t}_{uf} = \frac{\sum_{k=1}^{K} \hat{\alpha}_{f00}^{uk} - i_{zf} \sum_{k=1}^{K} a_{00}^k}{K} \qquad (21)$$

$$\hat{t}_{vf} = \frac{\sum_{k=1}^{K} \hat{\alpha}_{f00}^{vk} - j_{zf} \sum_{k=1}^{K} a_{00}^k}{K} \qquad (22)$$

Without loss of generality, we choose the object coordinate system in such a way that $\sum_{k=1}^{K} a_{00}^k = 0$ (the first frame only restricts the rotation and the component of the translation parallel to the image plane, so we have freedom to move the coordinate system along the axis orthogonal to the image plane). With this choice, we obtain the estimates

$$\hat{t}_{uf} = \frac{1}{K} \sum_{k=1}^{K} \hat{\alpha}_{f00}^{uk}, \qquad \hat{t}_{vf} = \frac{1}{K} \sum_{k=1}^{K} \hat{\alpha}_{f00}^{vk} \qquad (23)$$

**Optical flow parameters matrix** Replace the set of translation estimates $\{\hat{t}_{uf}, \hat{t}_{vf}\}$ given by (23) in the equation set (16-18). Define a set of parameters $\left\{\beta_f^{uk}, \beta_f^{vk}\right\}$ related to $\left\{\alpha_{f00}^{uk}, \alpha_{f00}^{vk}\right\}$ by

$$\beta_f^{uk} = \alpha_{f00}^{uk} - \frac{1}{K} \sum_{l=1}^{K} \alpha_{f00}^{ul}, \quad \beta_f^{vk} = \alpha_{f00}^{vk} - \frac{1}{K} \sum_{l=1}^{K} \alpha_{f00}^{vl} \qquad (24)$$

Collect the parameters $\left\{\beta_f^{uk}, \alpha_{fmn}^{uk}, \beta_f^{vk}, \alpha_{fmn}^{vk}\right\}$ in the matrix $R$, which we call the optical flow parameters matrix,

$$R = \begin{bmatrix} \beta_1^{u1} & \alpha_{110}^{u1} & \alpha_{101}^{u1} & \cdots & \beta_1^{uK} & \alpha_{110}^{uK} & \alpha_{101}^{uK} \\ \beta_2^{u1} & \alpha_{210}^{u1} & \alpha_{201}^{u1} & \cdots & \beta_2^{uK} & \alpha_{210}^{uK} & \alpha_{201}^{uK} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \beta_F^{u1} & \alpha_{F10}^{u1} & \alpha_{F01}^{u1} & \cdots & \beta_F^{uK} & \alpha_{F10}^{uK} & \alpha_{F01}^{uK} \\ \beta_1^{v1} & \alpha_{110}^{v1} & \alpha_{101}^{v1} & \cdots & \beta_1^{vK} & \alpha_{110}^{vK} & \alpha_{101}^{vK} \\ \beta_2^{v1} & \alpha_{210}^{v1} & \alpha_{201}^{v1} & \cdots & \beta_2^{vK} & \alpha_{210}^{vK} & \alpha_{201}^{vK} \\ \cdots & \cdots & \cdots & \cdots & \cdots & \cdots & \cdots \\ \beta_F^{v1} & \alpha_{F10}^{v1} & \alpha_{F01}^{v1} & \cdots & \beta_F^{vK} & \alpha_{F10}^{vK} & \alpha_{F01}^{vK} \end{bmatrix}$$

and the motion and shape parameters in matrices $M$ and $S$, as follows

$$M = \begin{bmatrix} i_{x1} & i_{x2} & \cdots & i_{xF} & j_{x1} & j_{x2} & \cdots & j_{xF} \\ i_{y1} & i_{y2} & \cdots & i_{yF} & j_{y1} & j_{y2} & \cdots & j_{yF} \\ i_{z1} & i_{z2} & \cdots & i_{zF} & j_{z1} & j_{z2} & \cdots & j_{zF} \end{bmatrix}^T$$

$$S^T = \begin{bmatrix} 0 & 1 & 0 & \cdots & 0 & 1 & 0 \\ 0 & 0 & 1 & \cdots & 0 & 0 & 1 \\ a_{00}^1 & a_{10}^1 & a_{01}^1 & \cdots & a_{00}^K & a_{10}^K & a_{01}^K \end{bmatrix} \qquad (27)$$

These enable us to factorize the optical flow matrix according to (16-18) and (24)

$$R = M S^T \qquad (28)$$

Matrix $R$ is $2F \times 3K$. In a noiseless situation, $R$ is rank deficient (rank 3) due to the redundancy of the observations imposed by the 3D rigid structure of the object. With noisy observations, $R$ may be full rank but the 3 largest eigenvalues should contain most of the energy.

**Factorization** To estimate the shape and motion parameters we use a two-step iterative method. The steps are: i) solve for shape with known positions, and ii) solve for positions with known shape. We will see that step i) leads to a linear LS problem and step ii), although nonlinear, decouples the estimation of the positions at different instants. The initialization is done by computing the *Singular Value Decomposition* (SVD) of the matrix $R$ and selecting the 3 largest eigenvalues. We get $R \simeq U \Sigma V^T$ where $U$ is $2F \times 3$, $\Sigma$ is diagonal $3 \times 3$, and $V^T$ is $3 \times 3K$. We initialize $M = U \Sigma^{\frac{1}{2}} A$ and $S^T = A^{-1} \Sigma^{\frac{1}{2}} V^T$ where $A$ is a non-singular $3 \times 3$ matrix determined by the constraints imposed by the structure of the matrix $M$ (expressions (9-14) and (26)). This step is similar to the procedure in reference [2], see [4] for the details.

**Shape estimate for known motion** The shape is described by the third row of the matrix $S^T$ which we denote by $s_3^T$. Given the matrix $M$, we denote by $\tilde{R}$ a matrix equal to $R$ at all entries except for the set $\left\{\alpha_{f10}^{uk}, \alpha_{f01}^{uk}, \alpha_{f10}^{vk}, \alpha_{f01}^{vk}\right\}$ which is replaced by $\left\{\gamma_{f10}^{uk}, \gamma_{f01}^{uk}, \gamma_{f10}^{vk}, \gamma_{f10}^{vk}\right\}$ defined as

$$\gamma_{f10}^{uk} = \alpha_{f10}^{uk} - i_{xf}, \qquad \gamma_{f01}^{uk} = \alpha_{f01}^{uk} - i_{yf} \qquad (29)$$

$$\gamma_{f10}^{vk} = \alpha_{f10}^{vk} - j_{xf}, \qquad \gamma_{f01}^{vk} = \alpha_{f01}^{vk} - j_{yf} \qquad (30)$$

According to expressions (25-30), we have

$$\tilde{R} = m_3 s_3^T \qquad (31)$$

where $m_3$ is the third column of the matrix $M$. The estimation of $s_3$ from $\tilde{R}$ and $m_3$, according to expression (31), is a linear problem. The LS solution is

$$\hat{s}_3^T = \left(m_3^T m_3\right)^{-1} m_3^T \tilde{R} = \frac{m_3^T \tilde{R}}{m_3^T m_3} \qquad (32)$$

**Motion estimate for known shape** From (25-27), note that the object position at instant $f$ only affects the rows $f$ and $f + F$ of matrix $R$. Given the object shape $S$, the LS estimate of the object rotation

$\{\theta_f, \phi_f, \psi_f\}$ for each frame $f$ is given by

$$c_f(\theta, \phi, \psi) = \left[ \begin{array}{c} r_f \\ r_{f+F} \end{array} \right] - \left[ \begin{array}{c} S[\ i_x \quad i_y \quad i_z\ ]^T \\ S[\ j_x \quad j_y \quad j_z\ ]^T \end{array} \right]$$

$$\left\{ \hat{\theta}_f\ \hat{\phi}_f, \hat{\psi}_f \right\} = \arg \min_{\{\theta, \phi, \psi\}} c_f^T(\theta, \phi, \psi) c_f(\theta, \phi, \psi) \quad (34)$$

To solve the non-linear minimization (34), we a Gauss-Newton method. Starting from an initial point $\{\theta_0, \phi_0, \psi_0\}$, the increments $\{\delta_\theta, \delta_\phi, \delta_\psi\}$ are found by minimizing (34) after truncating the Taylor series of $c_f(\theta_0 + \delta_\theta, \phi_0 + \delta_\phi, \psi_0 + \delta_\psi)$. Neglecting second and higher order terms, we have

$$c_f(\theta_0 + \delta_\theta, \phi_0 + \delta_\phi, \psi_0 + \delta_\psi) \simeq c_f(\theta_0, \phi_0, \psi_0) \quad (35)$$

$$+\nabla c_f(\theta_0, \phi_0, \psi_0) \left[ \begin{array}{c} \delta_\theta \\ \delta_\phi \\ \delta_\psi \end{array} \right]$$

where $\nabla$ is the gradient operator. Equating to zero the partial derivates of the cost function $c_f(\theta_0 + \delta_\theta, \phi_0 + \delta_\phi, \psi_0 + \delta_\psi)^T c_f(\theta_0 + \delta_\theta, \phi_0 + \delta_\phi, \psi_0 + \delta_\psi)$ with respect to the increments $\delta_\theta, \delta_\phi, \delta_\psi$, we obtain linear estimates $\hat{\delta}_\theta, \hat{\delta}_\phi, \hat{\delta}_\psi$ from the solution of

$$\nabla_{c_f}^T(\theta_0, \phi_0, \psi_0) \nabla c_f(\theta_0, \phi_0, \psi_0) \left[ \begin{array}{c} \hat{\delta}_\theta \\ \hat{\delta}_\phi \\ \hat{\delta}_\psi \end{array} \right] = \quad (36)$$

$$-\nabla_{c_f}^T(\theta_0, \phi_0, \psi_0) c_f(\theta_0, \phi_0, \psi_0)$$

When computing $\hat{\theta}_f, \hat{\phi}_f, \hat{\psi}_f$ we start the iterative process with the initial guess $\theta_0 = \hat{\theta}_{f-1}, \phi_0 = \hat{\phi}_{f-1}, \psi_0 = \hat{\psi}_{f-1}$. Initially, we have $\theta_1 = \phi_1 = \psi_1 = 0$ by definition. See [4] for the details.

## 4  Experiment

We used a real life video showing a static corner taped by a moving camera. Figure 2 shows two consecutive frames from the sequence of 10 images.
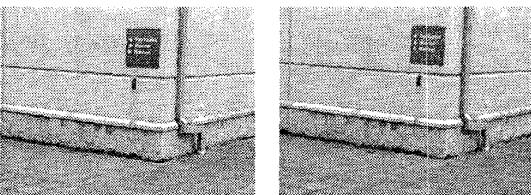


Figure 2: Image sequence.

Figure 3 shows a perspective view of the reconstructed 3D shape. It contains 3 planar patches: the floor and two walls. The angle between the walls is clearly seen. The angle of the walls with the floor can also be perceived.



Figure 3: Reconstructed 3D shape and texture.

## 5  Conclusions

We recover 3D structure from 2D video. The results obtained so far show that our method can be used in content-based video analysis tasks.

## References

[1] Radu S. Jasinschi and José M. F. Moura. Content-based video sequence representation. In *Proceedings of the IEEE International Conference on Image Processing*, Washington DC, U.S.A., October 1995.

[2] Carlo Tomasi and Takeo Kanade. Shape and motion from image streams under orthography: a factorization method. *International Journal of Computer Vision*, 9(2):137–154, 1992.

[3] Pedro M. Q. Aguiar and José M. F. Moura. Robust 3D structure from motion under ortography. In *Proceedings of the Tenth IEEE Image and Multidimensional Digital Signal Processing Workshop*, Alpbach, Austria, July 1998.

[4] Pedro M. Q. Aguiar and José M. F. Moura. 3D shape and motion recovery from induced parametric optical flow models. To be submitted.

[5] James R. Bergen, P. Anandan, Keith J. Hanna, and Rajesh Hingorani. Hierarchical model-based motion estimation. In *Proceedings of the Second European Conference on Computer Vision*, Santa Margherita Ligure, Italy, May 1992. Springer-Verlag.

[6] Serge Ayer. *Sequential and Competitive Methods for Estimation of Multiple Motions*. PhD thesis, École Polytechnique Fédérale de Lausanne, 1995.